

AD-A074 082

MASSACHUSETTS INST OF TECH LEXINGTON LINCOLN LAB

F/G 17/2

SPEECH ENHANCEMENT USING A SOFT-DECISION MAXIMUM LIKELIHOOD NOI--ETC(U)

JUN 79 R J MCAULAY, M L MALPASS

F19628-78-C-0002

UNCLASSIFIED

TN-1979-31

ESD-TR-79-163

NI

| OF |

AD
A074082



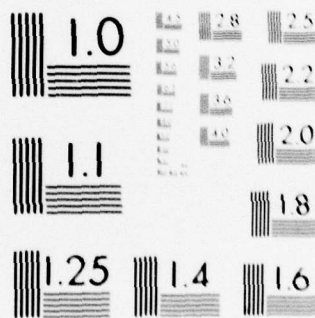
END

DATE

FILMED

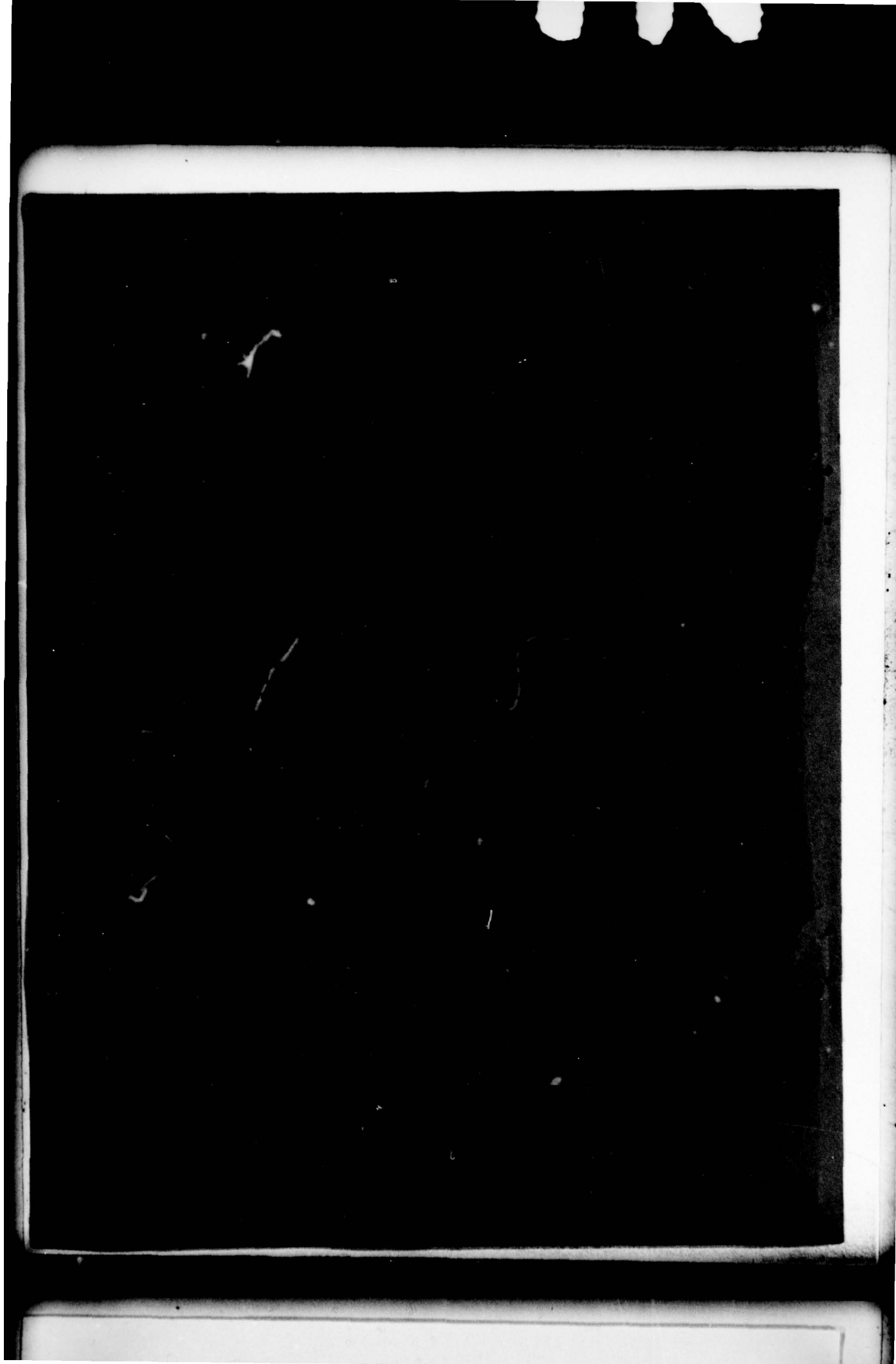
10-79

DDC



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

ADA074082



12

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
LINCOLN LABORATORY

SPEECH ENHANCEMENT USING A
SOFT-DECISION MAXIMUM LIKELIHOOD
NOISE SUPPRESSION FILTER

R. J. McAULAY
M. L. MALPASS

Group 24



TECHNICAL NOTE 1979-31

19 JUNE 1979

Approved for public release; distribution unlimited.

LEXINGTON

MASSACHUSETTS

ABSTRACT

One way of enhancing speech in an additive acoustic noise environment is to perform a spectral decomposition of a frame of noisy speech and to attenuate a particular spectral line depending on how much the measured speech plus noise power exceeds an estimate of the background noise. Using a two state model for the speech event (speech absent or speech present) and determining the maximum likelihood estimator of the speech power results in a new class of suppression curves which permits a tradeoff of noise suppression against speech distortion. The algorithm has been implemented in real time in the time domain, exploiting the structure of the channel vocoder. Extensive testing has shown that the noise can be made imperceptible by proper choice of the suppression factor.

Accession For	
NTIS GRA&I	<input checked="checked" type="checkbox"/>
DDC TAB	<input type="checkbox"/>
Unannounced	
Justification	
By _____	
Distribution/	
Availability Codes	
Dist.	Avail and/or special
A	

CONTENTS

ABSTRACT	iii
I. INTRODUCTION	1
II. ANALYSIS	3
A. Power Subtraction	5
B. Wiener Filtering	6
C. Maximum Likelihood Envelope Estimation	7
D. Two-State Soft Decision Maximum Likelihood Envelope Estimation	10
III. IMPLEMENTATION	16
IV. EXPERIMENTAL RESULTS AND CONCLUSIONS	25
ACKNOWLEDGMENT	28
APPENDIX: Modified Roberts Noise Detection Algorithm	29
REFERENCES	33

SPEECH ENHANCEMENT USING A SOFT-DECISION MAXIMUM LIKELIHOOD NOISE SUPPRESSION FILTER

I. INTRODUCTION

The need for secure military voice communication has led to the consideration of narrowband digital voice terminals. A preferred algorithm for this task is linear-predictive coding (LPC) which has demonstrated the ability to produce very intelligible speech with Diagnostic Rhyme Test (DRT) scores in excess of 90% at data rates as low as 2400 bps.[1] Unfortunately these results have been achieved only for clean speech, whereas many of the practical environments in which these terminals would be deployed, such as the airborne command post or the cockpits of jet fighter aircraft and helicopters, are characterized by a high ambient noise level, which in many cases causes the vocoded speech to suffer a significant degradation in intelligibility.[2] This has stimulated research into the problem of extracting the speech parameters (pitch, buzz-hiss and spectrum) from noisy speech in the hope that more robust algorithms could be found.[3,4,5]

Another approach to the noisy speech problem is to develop a prefilter that would enhance the speech prior to encoding so that the existing LPC vocoder could be applied in tandem without modification. Two general classes of algorithms have emerged: noise cancelling and noise suppression

prefilters. In the first case the coefficients of a tapped delay line are adapted to produce a minimum mean squared error estimate of the noise signal which is then subtracted from the noisy speech waveform to effect the noise cancellation.[6] In order to train the coefficients of the noise cancelling filter it is usually necessary to use a second microphone to provide a speech-free measurement of the background noise. Application of this technique to the cancellation of E4A advanced airborne command post noise has shown that although significant improvement in signal-to-noise ratio (SNR) can be obtained, the improvement in intelligibility, as measured by the Diagnostic Rhyme Test (DRT), is marginal.[7] Recent work by Sambur[8] has attempted to exploit the periodicity of voiced speech to eliminate the requirement for a second microphone. Thorough evaluation of this algorithm has not yet been published.

Considerably more work has been expended on the development of noise suppression prefilters. In this approach a spectral decomposition of a frame of noisy speech is performed and a particular spectral line is attenuated depending on how much the measured speech plus noise power exceeds an estimate of the background noise power.[9-13] Algorithms using the FFT have been tested against wideband noise and improvements in intelligibility have been indicated although no quantitative results have been given.[11] To date, the attenuation curves have been proposed on more or less an ad hoc basis, hence it is of interest to determine whether or not a more fundamental theoretical analysis could lead to a new suppression curve with substantially different properties. In the next section an

analytical model is proposed and used to determine the conditions under which the existing suppression curves can be justified. Having established a common basis, a new suppression curve is derived recognizing the fact that the degree of suppression should be weighted by the probability that a given measurement corresponds to speech plus noise or to noise alone. It is shown that a class of curves is obtained by varying the value of a suppression factor. This is a parameter that can be chosen to trade off noise suppression against speech distortion. The algorithm has been implemented in real time in the time domain, exploiting the structure of the channel vocoder to perform the spectral decomposition. Extensive testing has shown that the noise can be made imperceptible by proper choice of the suppression factor.

II. ANALYSIS

The prefilter design problem arises because a speech signal $s(t)$ has been corrupted by acoustically coupled background noise $w(t)$ to form the measurement $y(t) = s(t) + w(t)$. In speech it is not easy to specify a criterion which would lead to a "best" estimate of $s(t)$, hence a variety of algorithms are often proposed and evaluated by listening to the processed results. In order to provide a common theoretical basis for relating some of these algorithms it has been found useful to analyze the prefilter for a frame of data of length T ($T \sim 20$ millisec). A further simplification occurs by expanding $y(t)$ in terms of a set of basis functions $\{\phi_n(t)\}$ in such a way that the expansion coefficients are uncorrelated random variables. If

the covariance function of $y(t)$ is $R_y(t,u)$, then a suitable set of basis functions are obtained from the Karhunen-Loeve expansion,

$$\lambda(n)\phi_n(t) = \int_0^T R_y(t,u)\phi_n(u)du \quad 0 < t < T \quad (1)$$

Then on $(0,T)$

$$y(t) = \sum_{n=1}^{\infty} y_n \phi_n(t) \quad (2a)$$

$$y_n = \int_0^T y(t)\phi_n(t)dt = s_n + w_n \quad (2b)$$

Van Trees[14] shows that if the correlation time of $y(t)$ is less than the frame interval T , then an appropriate set of eigenfunctions and eigenvalues are

$$\phi_n(t) = \frac{1}{\sqrt{T}} \exp(j\frac{2\pi nt}{T}) \quad (3a)$$

$$\lambda(n) = S_y(\frac{n}{T}) \quad (3b)$$

where

$$S_y(f) = \int_0^T R_y(\tau)e^{-j2\pi f\tau} d\tau \quad (4)$$

is the power spectrum of the observed process. Since a narrowband vocoder

speech signals over a bandwidth less than a kilohertz. Finite number of expansion coefficients are used to characterize $y(t)$. The prefilter design problem then reduces to the problem of optimally extracting the speech signal $s(t)$ from the noisy observation $y = s + n$. If the speech and noise are modeled as independent Gaussian random processes, then the expansion coefficients are independent Gaussian random variables and we have:

$$x_n^2(t) = s_n^2(t) + n_n^2(t) \quad (5)$$

where

$$x_n(A) = \int_{-\infty}^{\infty} x_n(t) e^{-j\omega t} dt \quad (6a)$$

$$x_n(A) = \int_{-\infty}^{\infty} x_n(t) e^{-j\omega t} dt \quad (6b)$$

represent the power in the n th harmonic line of the speech and noise spectra.

A. POWER SUBTRACTION

Since it is well known that the perception of speech is phase insensitive, a reasonable criterion for a prefilter design is to produce the speech estimate

$$\hat{s}(t) = \sum_{n=1}^N \hat{s}_n \phi_n(t) \quad 0 \leq t \leq T \quad (7)$$

where $\hat{s}_n = \sqrt{\lambda_s(n)}$ since if $\lambda_s(n)$ were known, the spectrum of $\hat{s}(t)$ would be identical to the spectrum of $s(t)$. Of course, it is not known and provision must be made for estimating its value from an observation of y_n and knowledge of $\lambda_w(n)$. Since the probability density function for the complex Gaussian variate y_n is

$$p(y_n) = \frac{1}{\pi[\lambda_s(n) + \lambda_w(n)]} \exp \left\{ -\frac{|y_n|^2}{[\lambda_s(n) + \lambda_w(n)]} \right\} \quad (8)$$

then by maximizing $p(y_n)$ with respect to $\lambda_s(n)$ the maximum likelihood estimate of $\lambda_s(n)$ can be found to be

$$\hat{\lambda}_s(n) = |y_n|^2 - \lambda_w(n) \quad (9)$$

In order to maintain an identity system in the absence of noise, the input phase can be appended to the prefilter output by taking

$$\begin{aligned} \hat{s}_n &= \sqrt{\hat{\lambda}_s(n)} \frac{y_n}{|y_n|} \\ &= \left[\frac{|y_n|^2 - \lambda_w(n)}{|y_n|^2} \right]^{1/2} \cdot y_n \end{aligned} \quad (10)$$

which is known as the method of power subtraction. Modifications of this algorithm have been studied extensively by Boll [10], Preuss [12] and Berouti, et al [13].

B. Wiener Filtering

Whereas the power subtraction algorithm arises from an attempt to

obtain the best estimate of the speech spectrum, the Wiener Filter corresponds to the criterion of minimizing the mean squared error of best time domain fit to the speech waveform. Van Trees [15] has shown that this can be done by choosing the channel coefficients to be

$$\hat{s}_n = \frac{\lambda_s(n)}{\lambda_s(n) + \lambda_w(n)} \cdot y_n \quad (11)$$

Since the speech eigenvalues are unknown a priori, the maximum likelihood estimate developed in (8) can be used in (11) to result in the suppression rule

$$\hat{s}_n = \left[\frac{|y_n|^2 - \lambda_w(n)}{|y_n|^2} \right] \cdot y_n \quad (12)$$

which is simply the square of the suppression rule for the method of power subtraction.

C. Maximum Likelihood Envelope Estimation

The previous results were obtained assuming that the speech and the noise were independent Gaussian random processes. In the interest of exploring the importance of this assumption an alternative model is proposed in which the noise is a Gaussian random process while the speech is characterized by a deterministic waveform of unknown amplitude and phase. In this case the channel measurement is $y_n = s_n + w_n$ where now $s_n = A \exp(j\theta)$ where A determines the speech envelope and θ its phase. For the perception of speech an optimum estimate of its envelope is desired since

this would represent an estimate of the speech spectrum in the n th channel.

For Gaussian noise the probability density function of the channel measurement y_n is

$$p(y_n|A, \theta) = \frac{1}{\pi \lambda_w(n)} \exp \left[- \frac{|y_n|^2 - 2A \operatorname{Re}(e^{-j\theta} y_n) + A^2}{\lambda_w(n)} \right] \quad (13)$$

To obtain the maximum likelihood estimate of A , a maximum of $p(y_n|A, \theta)$ is sought. However the speech phase θ shows up as a nuisance parameter. Its effect can be eliminated by maximizing the average likelihood function

$$\overline{p(y_n|A)} = \int_0^{2\pi} p(y_n|A, \theta) p(\theta) d\theta \quad (14)$$

where $p(\theta)$ is the probability density function for the phase. Since it is reasonable to assume a uniform distribution on $(0, 2\pi)$, then the likelihood function for the spectral envelope becomes

$$\overline{p(y_n|A)} = \frac{1}{\pi \lambda_w(n)} \cdot \exp \left[- \frac{|y_n|^2 + A^2}{\lambda_w(n)} \right] \cdot \frac{1}{2\pi} \int_0^{2\pi} \exp \left[\frac{2A \operatorname{Re}(e^{-j\theta} y_n)}{\lambda_w(n)} \right] d\theta \quad (15)$$

The integral appearing in (15) is known as the modified Bessel function of the first kind and is labelled

$$I_0(|x|) = \frac{1}{2\pi} \int_0^{2\pi} \exp [\operatorname{Re}(e^{-j\theta} x)] d\theta \quad (16)$$

For large values of $|x|$ (≥ 3)

$$I_0(|x|) \sim \frac{1}{\sqrt{2\pi|x|}} \exp(|x|) \quad (17)$$

For this condition the likelihood function for the spectral envelope becomes

$$\overline{p(y_n|A)} = \frac{1}{\pi\lambda_w(n)} \cdot \frac{1}{\sqrt{2\pi \frac{2A|y_n|}{\lambda_w(n)}}} \cdot \exp \left[- \frac{|y_n|^2 - 2A|y_n| + A^2}{\lambda_w(n)} \right] \quad (18)$$

Maximizing this function with respect to A leads to the estimator

$$\hat{A} = \frac{1}{2} [|y_n| + \sqrt{|y_n|^2 - \lambda_w(n)}] \quad (19)$$

As before the input phase can be appended to this estimate of the envelope to produce the maximum likelihood estimate of the speech waveform

$$\begin{aligned} \hat{s}_n &= \hat{A} \frac{y_n}{|y_n|} \\ &= \left[\frac{1}{2} + \frac{1}{2} \sqrt{\frac{|y_n|^2 - \lambda_w(n)}{|y_n|^2}} \right] \cdot y_n \end{aligned} \quad (20)$$

D. Two-State Soft Decision Maximum Likelihood Envelope Estimation

The suppression rules for the power subtraction, Wiener filtering and maximum likelihood algorithms are illustrated in Fig. 1. Their suppression capabilities were evaluated for speech in airborne command post noise using a real time implementation of the prefilter (to be described in detail in Section III). While it was difficult to determine which algorithm did the best job of extracting the speech when speech was present, it was apparent that none of the algorithms adequately suppressed the background noise when speech was absent. This is hardly surprising in view of the fact that the suppression rules were derived on the assumption that speech was always present in the measured data. Had a detector been used to determine that a given frame of data consisted of noise alone, then obviously a better suppression rule would have been to apply greater attenuation than indicated by the curves in Fig. 1. From this point of view it follows that a better suppression curve might evolve if a two state model for the speech event is considered at the outset; that is either speech is present or it is not. Mathematically this leads to the binary hypothesis model:

$$\begin{aligned} H_0: \text{ speech absent: } |y_n| &= |w_n| \\ H_1: \text{ speech present: } |y_n| &= |Ae^{j\theta} + w_n| \end{aligned} \quad (21)$$

Only the measured envelope is used in this measurement model since it has already been shown that the measured phase provides no useful information in the suppression of the noise. A useful criterion for estimating the spectral

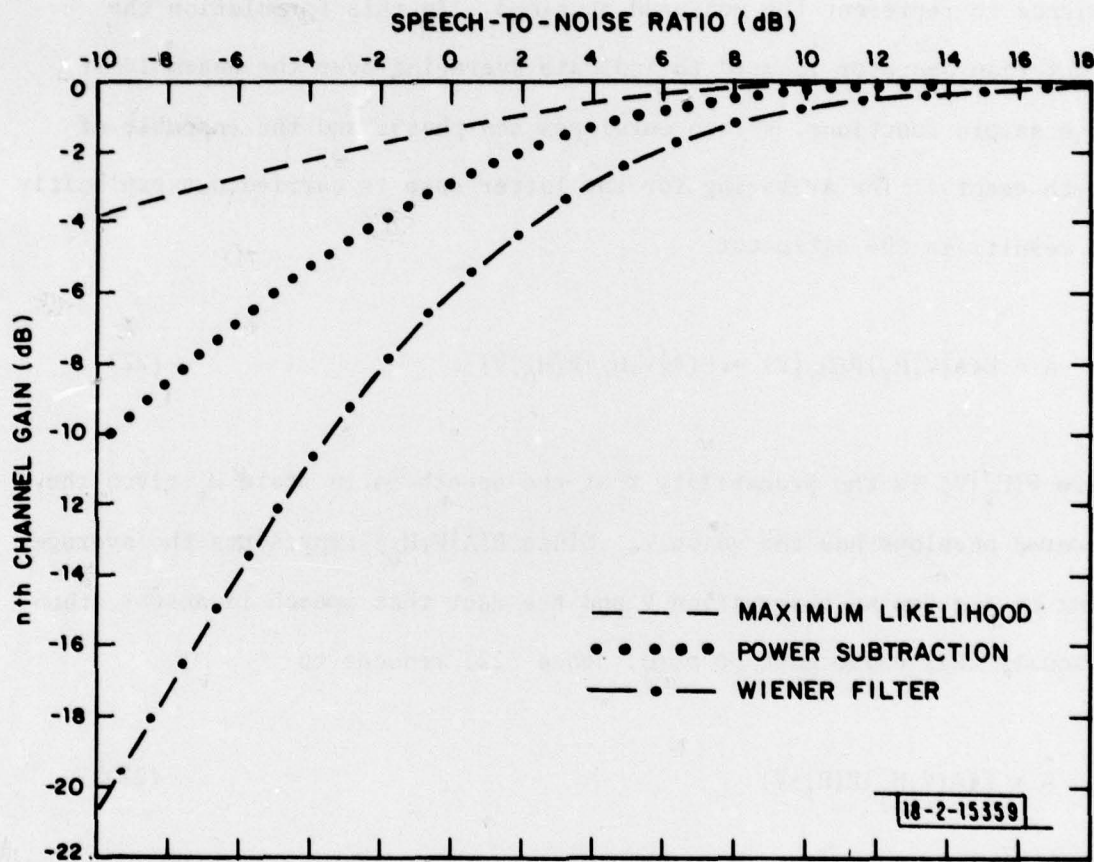


Fig. 1. Power subtraction, Wiener filter and maximum likelihood suppression rules.

envelope A , is to choose \hat{A} to minimize the mean squared spectral error $E(\hat{A} - A)^2$. It is well known [16] that the resulting estimator is the conditional mean $\hat{A} = E(A|V)$ where $V = |y_n|$ is used for notational convenience to represent the measured envelope. In this formulation the expectation operator is used to indicate averaging over the ensemble of noise sample functions, speech envelopes and phases and the ensemble of speech events. The averaging for the latter case is carried out explicitly and results in the estimator

$$\hat{A} = E(A|V, H_1)P(H_1|V) + E(A|V, H_0)P(H_0|V) \quad (22)$$

where $P(H_k|V)$ is the probability that the speech is in state H_k given the measured envelope has the value V . Since $E(A|V, H_0)$ represents the average value of A given an observation V and the fact that speech is absent, then obviously this value must be zero, hence (22) reduces to

$$\hat{A} = E(A|V, H_1)P(H_1|V) \quad (23)$$

Since $E(A|V, H_1)$ represents the minimum variance estimate of A when speech is present and since the maximum likelihood estimator is asymptotically efficient for large SNR, it suffices to replace $E(A|V, H_1)$ by the estimator derived in (19), hence

$$\hat{A} \sim \frac{1}{2} \left[V + \sqrt{V^2 - \lambda_w} \right] P(H_1|V) \quad (24)$$

Application of Bayes rule gives

$$P(H_1|V) = \frac{p(V|H_1)P(H_1)}{p(V|H_1)P(H_1) + p(V|H_0)P(H_0)} \quad (25)$$

where $p(V|H_k)$ is the a priori probability density function for the measured envelope given the speech state H_k . Assuming that the speech and noise states are equally likely (a worst case assumption),

$$P(H_1) = P(H_0) = \frac{1}{2} \quad (26)$$

Under hypothesis H_0 , $V = |w|$ and since the noise is complex Gaussian with mean zero and variance λ_w , it follows that the envelope has the Rayleigh pdf

$$p(V|H_0) = \frac{2V}{\lambda_w} \exp\left(-\frac{V^2}{\lambda_w}\right) \quad (27)$$

Under hypothesis H_1 , $V = |Ae^{j\theta} + w|$ and the envelope has the Rician pdf

$$p(V|H_1) = \frac{2V}{\lambda_w} \exp\left(-\frac{V^2 + A^2}{\lambda_w}\right) I_0\left(\frac{2AV}{\lambda_w}\right) \quad (28)$$

Defining the suppression factor ξ to be

$$\xi = \frac{A^2}{\lambda_w} \quad (29)$$

and substituting (26), (27), and (28) into (25) results in the following expression for the a posteriori probability for the presence of speech,

$$P(H_1|V) = \frac{\exp(-\xi) I_0 \left[2 \sqrt{\xi \left(\frac{V}{\lambda_w} \right)^2} \right]}{1 + \exp(-\xi) I_0 \left[2 \sqrt{\xi \left(\frac{V}{\lambda_w} \right)^2} \right]} \quad (30)$$

It is this term which contributes the soft suppression to the maximum likelihood envelope estimator. Appending the measured phase to the estimated envelope in order to preserve the identity system in the absence of noise, then the final suppression rule is

$$\begin{aligned} \hat{s} &= \hat{A} \frac{y}{|y|} \\ &= \left[\frac{1}{2} + \frac{1}{2} \sqrt{\frac{V^2 - \lambda_w}{V^2}} \right] \cdot P(H_1|V) \cdot y \end{aligned} \quad (31)$$

In Fig. 2 several curves for the a posteriori probability for the speech state $P(H_1|V)$ are illustrated for various values of the suppression factor ξ . The channel gains obtained when these a posteriori probabilities are

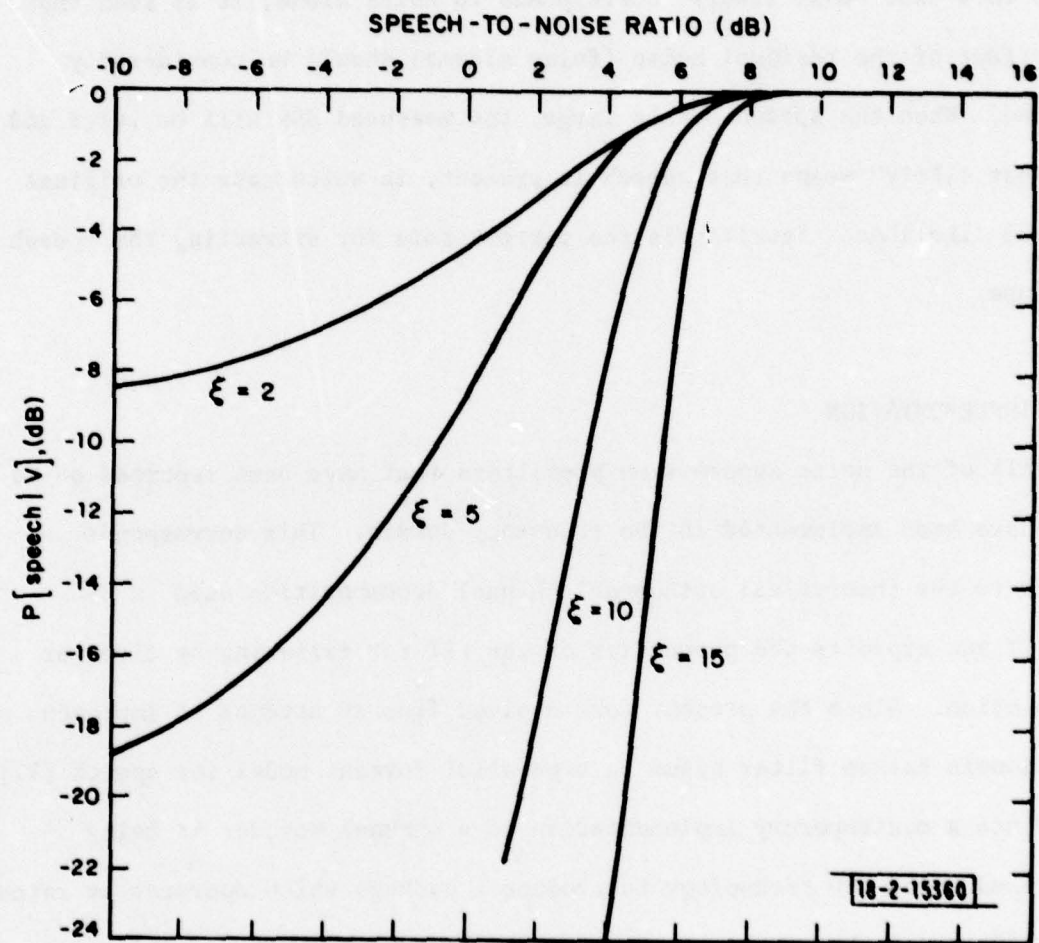


Fig. 2. A posteriori probability for the speech state.

appended to the maximum likelihood suppression rule are shown in Fig. 3. The two-state soft-decision maximum likelihood algorithm applies considerably more suppression when the measurement corresponds to low speech SNR. Since this case "most likely" corresponds to noise alone, it is seen that the effect of the residual noise (false alarms) should be considerably reduced. When the speech SNR is large, the measured SNR will be large and it "most likely" means that speech is present, in which case the original maximum likelihood algorithm is the correct rule for extracting the speech envelope.

III. IMPLEMENTATION

All of the noise suppression prefilters that have been reported on to date have been implemented in the frequency domain. This corresponds nicely to the theoretical orthogonal channel decomposition used in Section II and exploits the properties of the FFT for filtering by circular convolution. Since the present work evolved from an attempt to implement a time domain Kalman filter based on a parallel formant model for speech [17], and since a contemporary implementation of a channel vocoder is being developed using CCD technology to produce a package which operates at rates from 1.2 to 4.8 kbs, requires about 50 integrated circuits, occupies .22 cu. ft., requires 5 watts and weighs 5 lbs [18], it seemed appropriate to attempt a time domain implementation of the prefilter that could exploit this emerging technology. As in the channel vocoder 19 filters are used to span the frequency range 180 - 3720 Hz (the sampling rate was 7575 Hz).

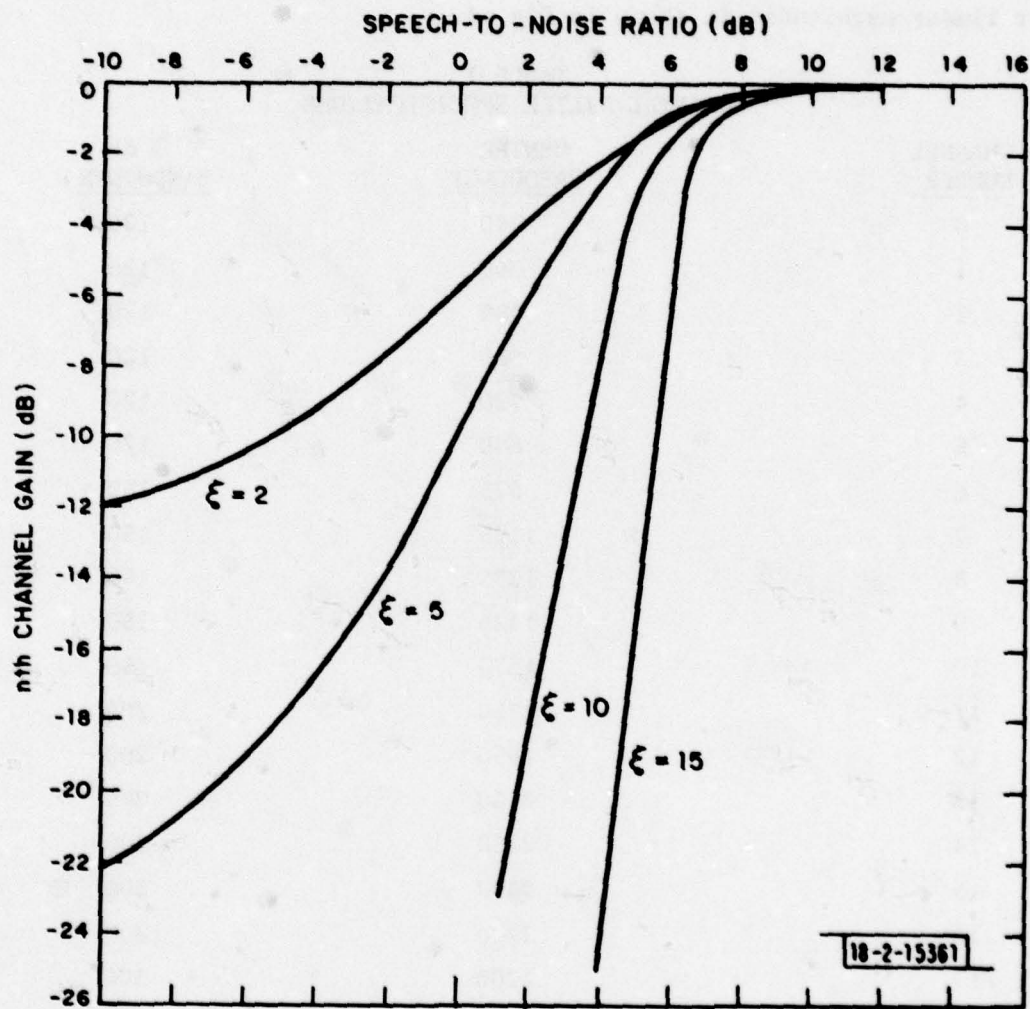


Fig. 3. Suppression rules for maximum likelihood with soft suppression.

Each filter in the bank is a result of a bandpass transformation of a second order Butterworth filter. The center frequencies and the bandwidths for each of the filters in the bank are listed in Table I and a plot of their linear magnitudes is shown in Fig. 4.

TABLE I
CHANNEL FILTER SPECIFICATIONS

<u>CHANNEL NUMBER</u>	<u>CENTER FREQUENCY</u>	<u>3 dB BANDWIDTH</u>
0	240	120
1	360	120
2	480	120
3	600	120
4	720	120
5	840	120
6	975	150
7	1125	150
8	1275	150
9	1425	150
10	1575	150
11	1750	200
12	1950	200
13	2150	200
14	2350	300
15	2600	300
16	2900	300
17	3200	300
18	3535	370

Sampling Rate = 132 μ sec

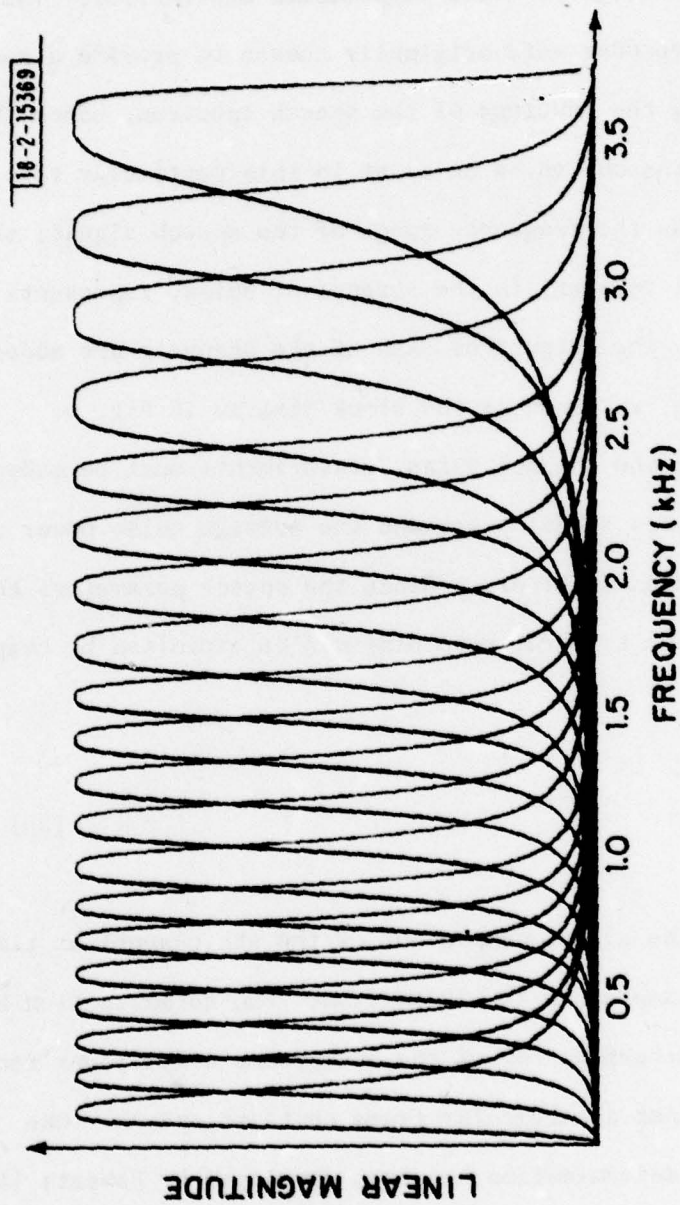


Fig. 4. The channel vocoder filter bank.

Although theory requires that the channels be orthogonal, in practice overlapping filters provide for spectral smoothing which is known to be an important factor in the design of noise suppression systems.[11] The filters in the channel vocoder were originally chosen to provide a good compromise for smoothing the envelope of the speech spectrum, hence their lack of orthogonality turns out to be an asset in this particular case. Since the 19 filters span the frequency range of the speech signal, the front end of the channel vocoder, in the absence of noise, represents an identity system provided the outputs of each of the channels are added alternately out of phase, as shown in the block diagram in Fig. 5.

In order to compute the channel gains, measurements must be made to determine the instantaneous signal power and the average noise power at the output of each of the channel filters. Since the speech parameters change very little in 20 ms, some temporal smoothing can be exploited by computing the signal power from

$$V^2 = \frac{1}{N} \sum_{k=1}^N y_n^2(k) \quad (32)$$

where $y_n(k)$ represents the signal sample out of the n th channel at time k , where there are N such samples in the 20 ms frame (the normalization by N will be unnecessary). Determination of the background noise power requires knowledge of whether or not a particular frame contains speech. One approach to making this determination has been developed by Roberts [19] who noted that a 4-sec histogram of the frame energies of the input signal

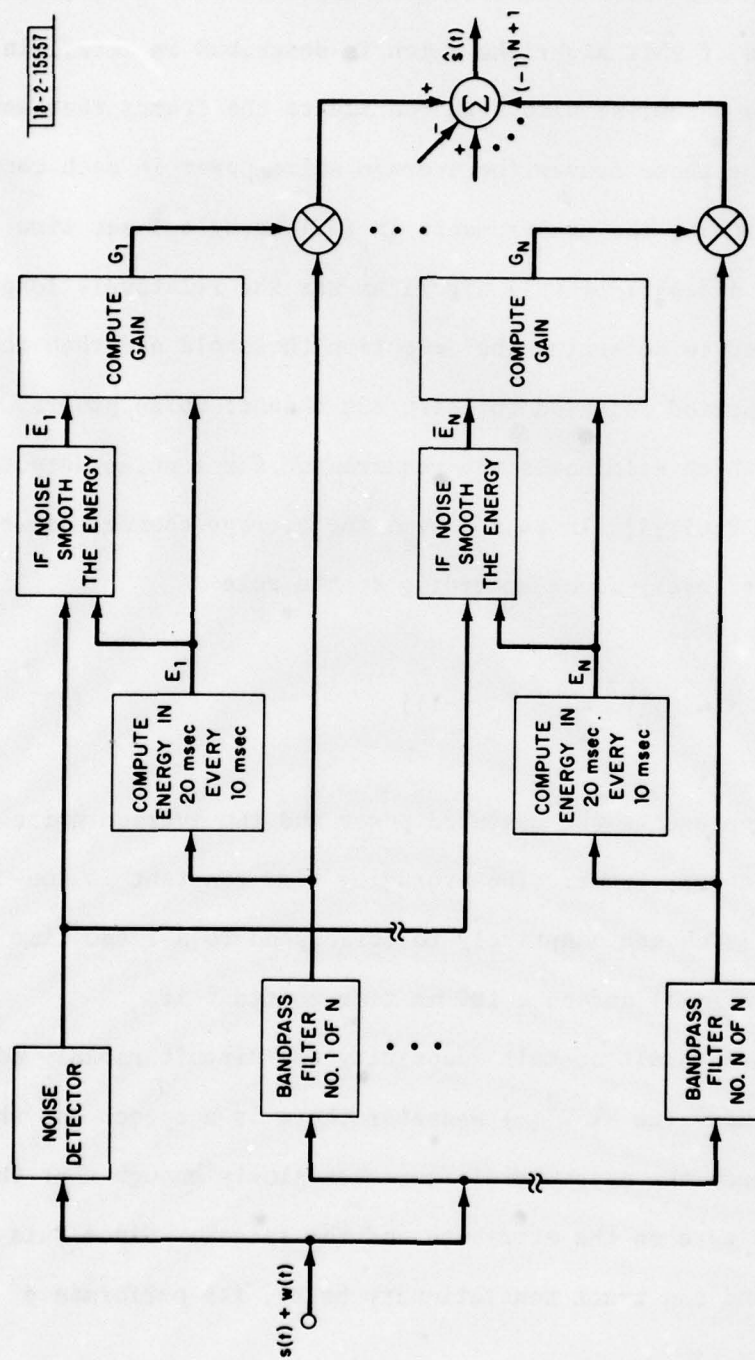


Fig. 5. Block diagram of the noise suppression prefilter.

was bimodal. He found that by setting a detection threshold between the modes, correct speech and noise classification could be made most of the time. A modification of this algorithm which is described in detail in the appendix, was used to determine with high confidence the frames that were absent of speech. For those frames the average noise power in each channel was estimated by smoothing the measurements in (32) using a 1-sec time constant. The major drawback of this algorithm was the relatively long adaptation time needed to determine the detection threshold and then the additional training period required to learn the channel noise powers. An alternative scheme, which eliminates the requirement for a noise detector has been proposed by Paul.[4] In this scheme the average channel noise power is updated after every frame according to the rule

$$\lambda_w(m) = \lambda_w(m-1) + \alpha(m)[V^2(m) - \lambda_w(m-1)] \quad (33)$$

where $V^2(m)$, $\lambda_w(m)$ represents the measured power and the average noise power computed for the mth frame. The averaging time constant is controlled by $\alpha(m)$ and is chosen adaptively to correspond to a 1 sec time constant if $V^2(m) \geq \lambda_w(m-1)$ and to a 100 ms time constant if $V^2(m) < \lambda_w(m-1)$. As a result of this adaptivity the circuit rapidly adapts the noise power to the value of $V^2(m)$ whenever there is a speech gap while during connected speech the noise level increases slowly enough that the noise power will not take on the attributes of the speech. Since this rule is easy to compute and can track nonstationary noise, its performance

warranted comparison with Roberts' noise detector.

Using the measurement of $V^2(m)$ and the estimated average value of $\lambda_w(m-1)$, the gain factor

$$g(m) = \frac{V^2(m) - \lambda_w(m-1)}{V^2(m)} \quad (34)$$

is computed for each channel.* Since $V^2(m)/\lambda_w(m-1)$ can be expressed in terms of $g(m)$ then the noise suppression rule, (30) and (31) can be written as

$$G = \frac{\hat{s}}{y} = \frac{1}{2}(1 + \sqrt{g}) \cdot \frac{\exp(-\xi) I_0\left(2\sqrt{\frac{\xi}{1-g}}\right)}{1 + \exp(-\xi) I_0\left(2\sqrt{\frac{\xi}{1-g}}\right)} \quad (35)$$

The advantage in using $g(m)$ as the independent variable is the fact that $0 \leq g(m) \leq 1$ which permits the use of a simple software divide routine in forming the normalization. For a given value of the suppression factor, ξ , the measured gain $g(m)$ is used as a pointer for a table look-up to determine the attenuation prescribed by (35). Fifteen tables corresponding to values of $\xi = 1, 2, 3, \dots, 15$ have been included in the prefilter with each table consisting of 50 values of the suppression rule computed for equal increments of $g(m)$ from 0 to 1. No attempt was made to optimize the design of these tables. All of the coding was done in machine language on the LDVT [19] which has the ability to key in a new value of the suppression

*This is where the normalization by N in (32) disappears.

factor in real time. This meant that the prefilter could easily be adjusted to accommodate a wide class of operational environments. This turned out to be a significant capability for effective noise suppression. Since the algorithm was designed to operate in real time, a 10 ms delay had to be incurred between the time the energies were measured and the time the corresponding gains could be computed and applied to the channel waveforms. This was done by computing the energies (block floating point) in 10 ms segments and adding consecutive segments together to produce the desired 20 ms energy measurement. This permitted computation of the raw gains, $G(m)$, every 10 ms. In order to avoid the introduction of discontinuities in the output waveform the final output is a smoothed gain $\bar{G}(m)$ obtained according to

$$\bar{G}(m) = \bar{G}(m-1) + \beta(m) [g(m) - \bar{G}(m-1)] \quad (36)$$

Since the introduction of smoothing can cause the prefilter to be slow to respond to a leading edge transition which could result in speech distortion, the gain in (35) is chosen adaptively according to the rule

$$\beta(m) = \begin{cases} 1 & \text{if } g(m) \geq \bar{G}(m-1) \\ \frac{1}{2} & \text{if } g(m) < \bar{G}(m-1) \end{cases} \quad (37)$$

In this way the prefilter responds immediately to an increase in the SNR which should minimize the potential for leading edge distortion. During a trailing edge, in which the gain will be decreasing, the smoothed gain will be used which will tend to maintain the speech signal even though the noise becomes dominant. It is the gain $\bar{G}(m)$ in (37) that is applied to the waveform at the output of each of the channel filters. These waveforms were then added together alternately 180° out of phase to produce the prefilter output waveform $\hat{s}(t)$.

IV. EXPERIMENTAL RESULTS AND CONCLUSIONS

Since the prefiltering algorithm operated in real time it was possible to perform extensive listening tests on a large speech and noise data base. It was of particular interest to determine the operational performance of the prefilter in conjunction with a 2400 bps vocoder operating in a background of E4A advanced airborne command post noise (ACPN). Source tapes were available for this environment consisting of lists spoken by six male speakers for which a DRT score and a diagnostic acceptability measure (DAM) could be computed. The recordings were made using both a high quality Altec microphone and a noise cancelling microphone.

The first experiment consisted of listening to the output of the prefilter for various values of the suppression factor. It was always possible to select a suppression factor which would render the background noise imperceptible, although, for cases in which the SNR was low enough, the cost in doing this was the introduction of various degrees of speech

distortion. In these cases, if the suppression factor was subsequently reduced, the speech distortion was reduced at the expense of introducing a perceptible level of background noise.

In the next experiment the prefilter was connected in tandem with the 2400 bps LPC vocoder which used the Gold-Rabiner pitch estimator.[21, 22] An unexpected result was obtained. If the suppression factor was set to remove the residual noise at the output of the prefilter then the speech quality at the vocoder output was poor due to both buzz-hiss errors and spectral distortion. If, however, the suppression factor was chosen so that the noise at the vocoder output was negligible, then a significantly lower value of the suppression factor was needed and the speech quality was quite good, although the Gold-Rabiner algorithm continued to make buzz-hiss errors, but at a lower rate. In other words, LPC itself has some suppression capabilities against weak noise which can usefully be exploited in the tandem connection. It was the flexibility in selecting the prefilter suppression factor which made this result possible.

Since the deployment of the LPC vocoder does allow for flexibility in the specification of the pitch extractor, it was of interest to determine whether or not algorithms that were specially designed to operate in noise would operate more effectively in the tandem connection. Such an algorithm, based on maximum likelihood estimation techniques, has been under development for some time [23] and was chosen to be tested against the Gold-Rabiner algorithm. In the subjective listening tests it was found that, indeed, smoother pitch tracks could be obtained with a lower rate of buzz-hiss errors.

Although the results of using the prefilter always produced subjectively more pleasant sounding speech to the ear since the annoying and tiresome background noise was removed, it was important to determine whether or not there was a corresponding quantitative improvement in intelligibility. To do this DRT scores are being obtained for the prefiltered speech and the speech out of the LPC tandem for both the Gold-Rabiner and the maximum likelihood pitch extraction algorithms. Results are currently being obtained for both the Altec dynamic microphone and the confidencer noise cancelling microphone and will be reported once all of the data has been collected and analyzed.

So far the focus has been on the 19-channel prefilter based on the principles of channel vocoder design. This was strictly a pragmatic choice which was made to facilitate the development of a real-time testbed. Questions relating to the number of filters, the bandwidths and the choice of center frequencies remain to be addressed. Although the time domain structure of the channel prefilter is well suited to an analog implementation using CCD technology, it is of interest to determine the tradeoffs with respect to a frequency domain approach using the FFT. Whatever candidate system is chosen for evaluation, using the class of suppression rules developed in this study allows the overall design to be optimized with respect to the noise suppression/speech distortion tradeoff by choosing an appropriate suppression factor. In this way performance differences can be attributed to the system design parameters independent of a particular suppression rule which may have represented a poor choice for the particular signal and noise conditions used in the evaluation.

ACKNOWLEDGMENT

The authors appreciate the technical interaction provided by their colleagues J. Tierney and T. Bially. Their intuition regarding the appropriateness of the channel vocoder front end as a preprocessing structure was particularly helpful. The authors would also like to thank A. J. McLaughlin for his support and encouragement throughout the course of this work.

APPENDIX

MODIFIED ROBERTS NOISE DETECTION ALGORITHM

In order to estimate the statistics of the background noise, it is desirable to inspect only those frames of data which have a high probability of containing no speech. To accomplish this, an adaptive energy threshold marking the probable boundary between noise and noise plus speech is established by monitoring the energy on a frame by frame basis and maintaining energy histograms which reflect the bimodal distribution of the energy. The flow chart for the algorithm, shown in Fig. 6, is described in the following paragraphs.

For each frame the sum of the squares of the input samples is computed. If this energy does not exceed 16 bits (i.e., does not strongly imply the presence of speech), the adaptive threshold algorithm is exercised. First, a decay factor of .995 is applied to a 128-bin histogram of uniform ranges of energy causing exponential decay of the histogram values with a time constant of 4 seconds. The value of the bin which encompasses the energy of the current frame is incremented by 160. A typical energy histogram after adaptation is complete is shown in Fig. 7a.

A second 128-point cumulative histogram is then formed to represent the area under the first histogram by computing the accumulated scores from the low energy bin through a high energy bin. Fig. 7b shows the result of accumulating the scores from the histogram in Fig. 7a. If the 10th point of the second histogram exceeds 25% of the total area, it is assumed that there is no noise present.

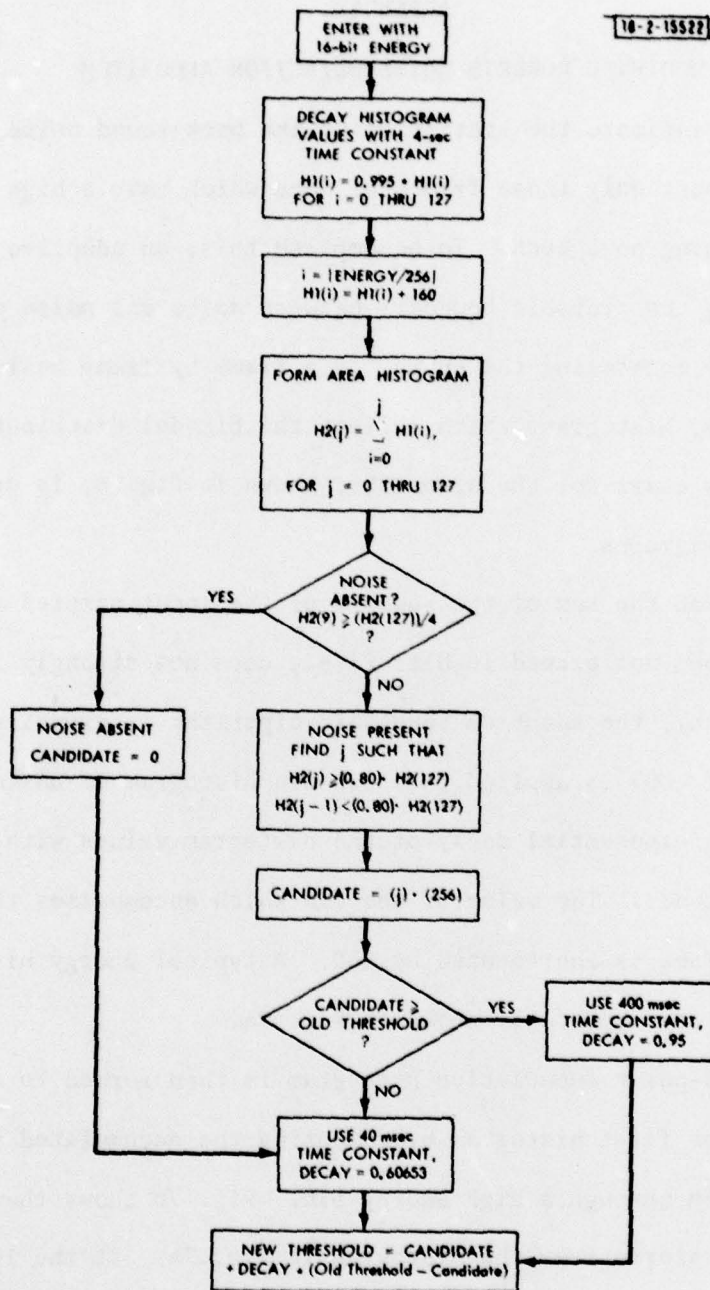


Fig. 6. Modified Roberts noise detection algorithm.

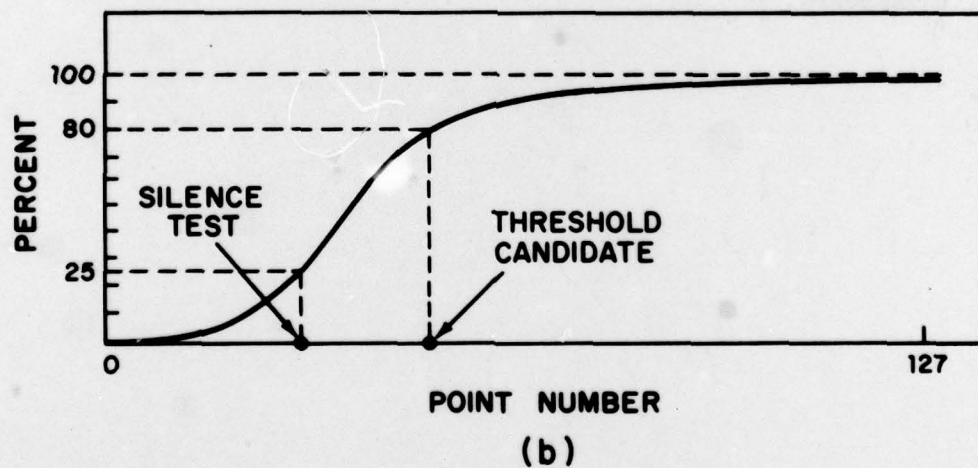
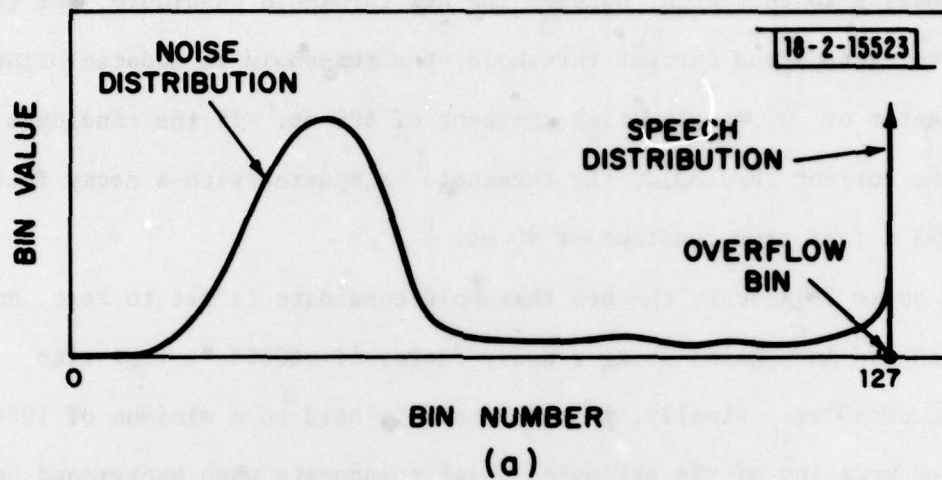


Fig. 7. (a) Typical energy histogram (H_1), (b) typical cumulative histogram (H_2).

If noise is present, a search is made through the second histogram for the point which represents 80% of the total area. The quantum of energy corresponding to this point becomes the new threshold candidate. If this candidate exceeds the current threshold, the threshold is updated using a decay factor of .95, a slow time constant of 400 ms. If the candidate is below the current threshold, the threshold is updated with a decay factor of .60653 a fast time constant of 40 ms.

If noise is absent, the new threshold candidate is set to zero, and the threshold is updated using a decay factor of .60653, a fast time constant of 40 ms. Finally, the threshold is held to a minimum of 1024 to guarantee updating of the estimated noise components when background noise suddenly disappears.

REFERENCES

1. T. E. Tremain, J. W. Fussell, R. A. Dean, B. M. Abzug, M. D. Cowing, and P. W. Boudra, Jr., "Implementation of Two Real Time Narrowband Speech Algorithms," EASCON '78 Record (25-27 September 1978) pp. 698-708.
2. C. Teacher and H. Watkins, "ANDVT Microphone and Audio System Study," Ketrion Final Report for the Commander, Naval Electronic System Command, Washington, DC (August 1978).
3. R. J. McAulay, "Maximum Likelihood Spectral Estimation Using State Variable Techniques," Proc. RADC Spectrum Estimation Workshop, (24-26 May 1978) pp. 63-68.
4. D. Paul, "A Robust Vocoder with Pitch-adaptive Spectral Envelope Estimation and an Integrated Maximum-Likelihood Pitch Estimator," Proc. International Conference on Acoustics, Speech, and Signal Processing (2-4 April 1979) pp. 64-68.
5. J. S. Lim and A. V. Oppenheim, "All-Pole Modelling of Degraded Speech," IEEE Trans. Acoust., Speech and Signal Processing ASSP-26, 197 (1978).
6. B. Widrow, et al, "Adaptive Noise Cancelling: Principles and Applications," Proc. IEEE 63, 1692 (1978).
7. R. C. Rohlf, "Speech Enhancement Using the Widrow-Hoff Adaptive Noise-cancelling Tapped Delay-line Filter; a CSP30 Implementation," Mitre Technical Report, MTR-3626, (March 1979).
8. M. R. Sambur, "Adaptive Noise Cancelling for Speech Signals," IEEE Trans. Acoust., Speech and Signal Processing ASSP-26, 419 (1978).
9. M. R. Weiss, E. Aschkenasy and T. W. Parsons, "Study and Development of the INTEL Technique for Improving Speech Intelligibility," Technical Report No. RADC-TR-75-108, RADC, Griffiss Air Force Base New York (April 1975).
10. S. F. Boll, "Suppression of Acoustic Noise in Speech Using Spectral Subtraction," IEEE Trans. Acoust., Speech and Signal Processing ASSP-26, 113 (1978).

11. R. A. Curtis and R. J. Niederjohn, "An Investigation of Several Frequency-domain Processing Methods for Enhancing the Intelligibility of Speech in Wideband Random Noise," Proceedings of International Conference on Acoustics, Speech and Signal Processing, (10-12 April 1978) pp. 602-605.
12. R. D. Preuss, "A Frequency Domain Noise Cancelling Preprocessor for Narrowband Speech Communications Systems," Proc. International Conference on Acoustics, Speech and Signal Processing, (2-4 April 1979) pp. 212-215.
13. M. Berouti, R. Schwartz and J. Makhoul, "Enhancement of Speech Corrupted by Acoustic Noise," Proc. International Conference on Acoustics, Speech and Signal Processing, (2-4 April 1979), pp. 208-211.
14. H. L. Van Trees, Detection, Estimation and Modulation Theory, Part I, (Wiley, New York, 1968), pp. 205-207.
15. *ibid.*, pp. 198-206.
16. *op cit.*, pp. 54-56.
17. R. J. McAulay, "A Structure for Robust Speech Processing," Proc. International Conference on Communications, Vol. 1 (4-7 June 1978) pp. 8.5.1-8.5.3.
18. P. E. Blankenship, "An NMOS LSI Channel Vocoder Implementation," Proc. EASCON '78, 25-27 September 1978, pp. 684-692, and SPIE Technical Symposium East '79, April 17-20, 1979, Washington, DC.
19. J. Roberts, "Modification to Piecewise LPC," MITRE Working Paper, WP-21752 (12 May 1978).
20. P. E. Blankenship, "LDVT: High Performance Minicomputer for Real-time Speech Processing," EASCON '77 Record, 26 September 1977.
21. B. Gold and L. Rabiner, "Parallel Processing Techniques for Estimating Pitch Periods of Speech in the Time Domain," J. Acoust. Soc. Am., 46, 442 (1969).
22. E. M. Hofstetter, P. E. Blankenship, M. L. Malpass and S. Seneff, "Vocoder Implementations on the Lincoln Digital Voice Terminal," EASCON '77 Record, (26 September 1977).
23. R. J. McAulay, "Design of a Robust Maximum Likelihood Pitch Estimator for Speech in Additive Noise," Technical Note 1979-28, Lincoln Laboratory, M.I.T. (11 June 1979).

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

18 19 REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER ESD-TR-79-163	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) Speech Enhancement Using a Soft-Decision Maximum Likelihood Noise Suppression Filter	5. TYPE OF REPORT & PERIOD COVERED Technical Note	6. PERFORMING ORG. REPORT NUMBER Technical Note 1979-31
7. AUTHOR(s) Robert J. McAulay and Marilyn L. Malpass	8. CONTRACT OR GRANT NUMBER(s) F19628-78-C-0002	
9. PERFORMING ORGANIZATION NAME AND ADDRESS Lincoln Laboratory, M.I.T. P.O. Box 73 Lexington, MA 02173	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS Program Element No. 33401F Project No. 7280	
11. CONTROLLING OFFICE NAME AND ADDRESS Air Force Systems Command, USAF Andrews AFB Washington, DC 20331	12. REPORT DATE 19 Jun 1979	13. NUMBER OF PAGES 40
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) Electronic Systems Division Hanscom AFB Bedford, MA 01731	15. SECURITY CLASS. (of this report) Unclassified	15a. DECLASSIFICATION DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES None		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) narrowband digital speech LPC noise suppression speech enhancement prefilter		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) One way of enhancing speech in an additive acoustic noise environment is to perform a spectral decomposition of a frame of noisy speech and to attenuate a particular spectral line depending on how much the measured speech plus noise power exceeds an estimate of the background noise. Using a two state model for the speech event (speech absent or speech present) and determining the maximum likelihood estimator of the speech power results in a new class of suppression curves which permits a tradeoff of noise suppression against speech distortion. The algorithm has been implemented in real time in the time domain, exploiting the structure of the channel vocoder. Extensive testing has shown that the noise can be made imperceptible by proper choice of the suppression factor.		

207650

y/B